

REC'D 07 JAN 2005

WIPO

PCT

7B/2005/050038

PA 1125716

THE UNITED STATES OF AMERICA

TO ALL TO WHOM THESE PRESENTS SHALL COME:

UNITED STATES DEPARTMENT OF COMMERCE

United States Patent and Trademark Office

February 09, 2004

THIS IS TO CERTIFY THAT ANNEXED HERETO IS A TRUE COPY FROM THE RECORDS OF THE UNITED STATES PATENT AND TRADEMARK OFFICE OF THOSE PAPERS OF THE BELOW IDENTIFIED PATENT APPLICATION THAT MET THE REQUIREMENTS TO BE GRANTED A FILING DATE UNDER 35 USC 111.

APPLICATION NUMBER: 60/535,113 ✓

FILING DATE: January 08, 2004 ✓

PRIORITY DOCUMENT

SUBMITTED OR TRANSMITTED IN
COMPLIANCE WITH RULE 17.1(a) OR (b)



By Authority of the
COMMISSIONER OF PATENTS AND TRADEMARKS

T. Wallace
T. WALLACE
Certifying Officer

BEST AVAILABLE COPY

PATENT APPLICATION SERIAL NO. _____

U.S. DEPARTMENT OF COMMERCE
PATENT AND TRADEMARK OFFICE
FEE RECORD SHEET

01/12/2004 KBETEM1 00000085 141270 60535113

01 FC:1005 160.00 DA

PTO-1556
(5/87)

010804

17169

Please type a plus sign (+) inside this box → ☐ +

Approved for use through 10/31/2002. OMB 0851-0032
Patent and Trademark Office; U.S. DEPARTMENT OF COMMERCE

Under the Paperwork Reduction Act of 1995, no persons are required to respond to a collection of information unless it displays a valid OMB control number.

PROVISIONAL APPLICATION FOR PATENT COVER SHEET

This is a request for filing a PROVISIONAL APPLICATION FOR PATENT under 37 CFR 1.53 (c).

Express Mail Label No. EV 312 069 508		Date of Deposit: January 8, 2004	
INVENTOR(S)			
Given Name (first and middle [if any])	Family Name or Surname	Residence (City and either State or Foreign Country)	
Reinder Jaap Dietwig Jos Clement	BRIL LOWET	WAALRE, NETHERLANDS EINDHOVEN, NETHERLANDS	
<input type="checkbox"/> Additional inventors are being named on the _____ separately numbered sheets attached hereto			
TITLE OF THE INVENTION (280 characters max)			
RESOURCE-MANAGEMENT IN A MULTI-PROCESSOR SYSTEM			
CORRESPONDENCE ADDRESS			
Direct all correspondence to:			
<input checked="" type="checkbox"/> Customer Number	24737	*24737*	
OR	Type Customer Number here		
<input checked="" type="checkbox"/> Firm or Individual Name	U. S. PHILIPS CORPORATION		
Address	P. O. BOX 3001		
Address			
City	BRIARCLIFF MANOR	State	NY
Country	USA	ZIP	10510
Telephone	(914) 945-6000	Fax	(914) 332-0615
ENCLOSED APPLICATION PARTS (check all that apply)			
<input checked="" type="checkbox"/> Specification Number of Pages	24	<input type="checkbox"/> CD(s), Number	
<input checked="" type="checkbox"/> Drawing(s) Number of Sheets	6	Other (specify)	
<input type="checkbox"/> Application Data Sheet. See 37 CFR 1.76			
METHOD OF PAYMENT OF FILING FEES FOR THIS PROVISIONAL APPLICATION FOR PATENT (check one)			
<input type="checkbox"/> Applicant claims small entity status. See 37 CFR 1.27.			
<input type="checkbox"/> A check or money order is enclosed to cover the filing fees			
<input checked="" type="checkbox"/> The Commissioner is hereby authorized to charge filing fees or credit any overpayment to Deposit Account Number:	14-1270	FILING FEE AMOUNT (\$) 160.00	
<input type="checkbox"/> Payment by credit card. Form PTO-2038 is attached.			
The invention was made by an agency of the United States Government or under a contract with an agency of the United States Government.			
<input checked="" type="checkbox"/> No.			
<input type="checkbox"/> Yes, the name of the U.S. Government agency and the Government contract number are: _____			

Respectfully submitted,
SIGNATURE

Date January 8, 2004

TYPED or PRINTED NAME DANIEL J. PIOTROWSKI

REGISTRATION NO.: 42,079
(if appropriate)

TELEPHONE (914) 333-9624

Docket Number: US 040013

USE ONLY FOR FILING A PROVISIONAL APPLICATION FOR PATENT

This collection of information is required by 37 CFR 1.51. The information is used by the public to file (and by the PTO to process) a provisional application. Confidentiality is governed by 35 U.S.C. 122 and 37 CFR 1.14. This collection is estimated to take 8 hours to complete, including gathering, preparing, and submitting the complete provisional application to the PTO. Time will vary depending upon the individual case. Any comments on the amount of time you require to complete this form and/or suggestions for reducing this burden, should be sent to the Chief Information Officer, U.S. Patent and Trademark Office, U.S. Department of Commerce, Washington, D.C., 20231. DO NOT SEND FEES OR COMPLETED FORMS TO THIS ADDRESS. SEND TO: Box Provisional Application, Assistant Commissioner for Patents, Alexandria, VA 22313.

010804

RESOURCE MANAGEMENT IN A MULTI-PROCESSOR SYSTEM

5 The present invention relates to a resource management method and apparatus that is particularly, but not exclusively, suited to resource management of multi-processor real-time systems.

 For systems-on-silicon (SoS) or systems-on-chip (SoC), memory is becoming a dominant limiting factor, since, from the point of view of the amount of silicon needed, adding another processing core is no longer a problem. As a consequence, a SoS or SoC
10 may contain multiple processor cores.

 The management of memory is a crucial aspect of resource management for a multiprocessor system. Various methods have been developed to optimize memory use for single processor systems, including generalizing the use of preemption points to the
15 management of main memory, especially in real-time systems. In these approaches, rather than preempting tasks at arbitrary moments during their execution, those tasks are preferably only preempted at dedicated preemption points based on their memory usage.

 In a typical prior art approach, suspension of a task is referred to as task preemption, or preemption of a task, and the term "task" is used to denote a unit of
20 execution that can compete on its own for system resources such as memory, CPU, I/O devices, etc. For purposes of discussion, a task is assumed to be a succession of continually executing jobs, each of which comprises one or more sub-jobs. For example, a task can comprise "demultiplexing a video stream", and involve reading-in incoming streams, processing the streams and outputting corresponding data. These steps are carried
25 out with respect to each incoming data stream, so that reading, processing and outputting with respect to a single stream corresponds to performing one job, each having three sub-jobs. Thus, when there is a plurality of packets of data to be read-in and processed, the job would be performed a corresponding plurality of times. A sub-job can be considered to relate to a functional component of the job and in a multiprocessor system each stream
30 would be assigned to a different processor or subset of processors.

 A known method of scheduling a plurality of tasks in a data processing system requires that each sub-job of a task have a set of suspension criteria, called suspension data, that specifies the processing preemption points and corresponding conditions for

suspension of a sub-job based on its memory usage [4] [5]. The amount of memory that is used by the data processing system is thus indirectly controlled by this suspension data, via these preemption points, which specify the amounts of memory required at these preemption points in a job's execution.

5 Thus, these preemption points can be utilized to avoid data processing system crashes due to a lack of memory. When a real-time task is characterized as comprising a plurality of sub-jobs, its preemption points preferably coincide with the sub-job boundaries of the task.

10 Data indicative of memory usage of a task conforming to the suspension data associated with each sub-job of a task can, for example, be embedded into a task via a line of code that requests a descheduling event, specifying that a preemption point has been reached in the processing of the task, i.e., a sub-job boundary has been reached. That is, the set of start points of the sub-jobs of a task constitute a set of preemption points of that task. The j^{th} preemption point $P_{i,j}$ of a task τ_i is characterized by information related to the
15 preemption point itself and information related to the succeeding non-preemptible sub-job interval $I_{i,j}$ between the j^{th} preemption point and the next preemption point, i.e., the $(j+1)^{\text{th}}$ preemption point.

20 At run time, a task informs the controlling operating system when it arrives at preemption points, e.g. when it starts a sub-job, switches between sub-jobs, and completes a sub-job, and the operating system decides when and where execution of a task is preempted. Ideally, preemption may occur at a preemption point or at any other point during the execution of a task. However, such flexibility of choice of preemption comes at the cost of consistency so that preemption is limited to preemption points to maintain consistency.

25 A prior art preemption point approach based on main memory requirements that does not jeopardize consistency of the system, necessarily limits the preemption of all tasks to their preemption points and matching synchronization primitives for controlling exclusive use of a resource to both be within a sub-job boundary. As is known in the art, a component (e.g. a software component, which can comprise one or more tasks) can have a
30 programmable interface that comprises the properties, functions or methods and events that the component defines [6]. For purposes of discussion, a task is assumed to be accompanied by an interface 100 that includes, at a minimum, main memory data required by the task, $MP_{i,j}$ 101b, as illustrated in FIG. 1.

For the purposes of discussion, a task is assumed to be periodic and real-time, and characterized by a period T and a phasing F , where $0 \leq F < T$, which means that a task comprises a sequence of sub-jobs, each of which is released at time $F+nT$, where $n = 0 \dots N$. As an example only and as illustrated in FIG. 2, set-top box 200 is assumed to execute three tasks – (1) display menu on the User Interface 205, (2) retrieve text information from a content provider 203, and (3) process some video signals – and each of these 3 tasks is assumed to comprise a plurality of sub-jobs. For ease of presentation, it is assumed that the sub-jobs are executed sequentially.

At least some of these sub-jobs can be preempted and the boundaries between these sub-jobs that can be preempted provide preemption points and are summarized in Table 1:

Task τ_i	Task description	Number of preemptable sub-jobs of task τ_i $m(i)$
τ_1	display menu on the GUI	3
τ_2	retrieve text information from content provider	2
τ_3	process video signals	2

TABLE 1

Referring also to FIG.3, for each task the suspension data 101 comprises: information relating to a preemption-point $P_{i,j}$ 301, such as the maximum amount of memory $MP_{i,j}$ 302 required at the preemption point, and information relating to the interval $I_{i,j}$ 303 between successive preemption-points, such as the worst-case amount of memory $MI_{i,j}$ 304 required in an intra-preemption point interval (i represents task τ_i and j represents a preemption point).

More specifically, suspension data 101 comprises data specifying

1. preemption point j of the task τ_i ($P_{i,j}$) 101a;
2. maximum memory requirements of task τ_i , $MP_{i,j}$, at preemption point j of that task, where $1 \leq j \leq m(i)$ 101b;

3. interval, $I_{i,j}$, between successive preemption points j and $(j+1)$ corresponding to sub-job j of task τ_i , where $1 \leq j \leq m(i)$ 101c; and
4. maximum (i.e. worst-case) memory requirements of task τ_i , $MI_{i,j}$, in the interval j of that task, where $1 \leq j \leq m(i)$ 101d.

5 Table 2 illustrates the suspension data 101 for the current example (each task has its own interface, so that in the current example, the suspension data 101 corresponding to the first task τ_1 comprises the data in the first row of Table 2, the suspension data 101 corresponding to the second task τ_2 comprises the second row of Table 2, etc.):

Task τ_i	$MP_{i,1}$	$MI_{i,1}$	$MP_{i,2}$	$MI_{i,2}$	$MP_{i,3}$	$MI_{i,3}$
τ_1	0.2	0.6	0.2	0.4	0.1	0.6
τ_2	0.1	0.5	0.2	0.8	-	-
τ_3	0.1	0.2	0.1	0.3	-	-

TABLE 2

Suppose a set-top box 200 is equipped with 1.5 Mbytes of memory. Under normal, or non-memory based preemption conditions, this set-top box 200 behaves as follows.

15 Referring now to FIG. 4A, a processor 401 may be expected to schedule tasks according to some sort of time slicing or priority based preemption, meaning that all 3 tasks run concurrently, i.e. effectively at the same time. It is therefore possible that each task can be scheduled to run its most memory intensive sub-job at the same time. The worst-case memory requirement of these three tasks, M^P , is given by:

20

$$M^P = \sum_{i=1}^3 \max_{j=1}^{m(i)} MI_{i,j} \quad (\text{Equation 1})$$

For tasks τ_1 , τ_2 and τ_3 M^P is thus the maximum memory requirements of τ_1 (being $MI_{1,1}$) plus the maximal memory requirements of task τ_2 (being $MI_{2,2}$) plus the maximal memory requirements of task τ_3 (being $MI_{3,2}$). These maximum requirements are indicated by the Table 2 entries in bold:

$$M^P = 0.6 + 0.8 + 0.3 = 1.7 \text{ Mbytes.}$$

This exceeds the memory available to the set-top box 200 by 0.2 Mbytes, so that, in the absence of any precautionary measures and if these sub-jobs are to be processed at the same time, the set-top box 200 crashes.

5 Referring now to FIG. 5, suppose tasks are scheduled in accordance with a scheduling algorithm and a data structure is maintained for each task τ_i after it has been created and that matching pairs of primitives for exclusive use of a resource do not span a sub-job boundary. Suppose further that a scheduler 501 employs a conventional priority-based, preemptive scheduling algorithm, which essentially ensures that, at any point in
10 time, the currently running task is the one with the highest priority among all ready-to-run tasks in the system. As is known in the art, the scheduling behavior can be modified by selectively enabling and disabling preemption for the running, or ready-to-run, tasks.

A task manager 503 receives the suspension data 101 corresponding to a newly received task and evaluates whether preemption is required or not and if it is required,
15 passes this newly received information to the scheduler 501, requesting preemption. Suppose details of the tasks are as defined in Table 2, and assume that task τ_1 (and only τ_1) is currently being processed and that the scheduler is initially operating in a mode in which there are no memory-based constraints.

Suppose now that task τ_2 is received by the task manager 503, which reads the
20 suspension data 101 from its interface Int₂ 100, and identifies whether or not the scheduler 501 is working in accordance with memory-based preemption. Since, in this example, it is not, the task manager 503 evaluates whether the scheduler 501 needs to change to memory-based preemption. This therefore involves the task manager 503 retrieving worst case suspension data corresponding to all currently executing tasks (in this example task τ_1)
25 from a suspension data store 505, evaluating Equation 1 and comparing the evaluated worst-case memory requirements with the memory resources available. Continuing with the example introduced in Table 2, Equation 1, for τ_1 and τ_2 , is:

$$M^P = \sum_{i=1}^2 \max_{j=1}^{m(i)} MI_{i,j} = 0.6 + 0.8 = 1.4 \text{ Mbytes}$$

30

This is less than the available memory, so there is no need to change the mode of operation of the scheduler 501 to memory-based preemption (i.e. there is no need to

constrain the scheduler based on memory usage). Thus, if the scheduler 501 were to switch between task τ_1 and task τ_2 – e.g. to satisfy execution time constraints of task τ_2 , meaning that both tasks effectively reside in memory at the same time – the processor never accesses more memory than is available.

5 Next, and before tasks τ_1 and τ_2 have completed, another task τ_3 is received. The task manager 503 reads the suspension data 101 from interface Int_3 associated with the task τ_3 , evaluating whether the scheduler 501 needs to change to memory-based preemption. Assuming that the scheduler 501 is multi-tasking tasks τ_1 and τ_2 , the worst case memory requirements for all three tasks is now

10

$$M^P = \sum_{i=1}^3 \max_{j=1}^{m(i)} MI_{i,j} = 0.6 + 0.8 + 0.3 = 1.7 \text{ Mbytes}$$

This exceeds the available memory, so the task manager 503 requests and retrieves memory usage data $MP_{i,j}, MI_{i,j}$ 101b, 101d for all three tasks from the suspension data store 15 505, and evaluates whether, based on this retrieved memory usage data, there are sufficient memory resources to execute all three tasks. This can be ascertained through evaluation of the following equation:

$$\begin{aligned} M^D &= \sum_{i=1}^3 \max_{j=1}^{m(i)} MP_{i,j} + \max_{i=1}^3 (\max_{j=1}^{m(i)} MI_{i,j} - \max_{j=1}^{m(i)} MP_{i,j}) \quad (\text{Equation 2}) \\ &= 0.2 + 0.2 + 0.1 + \max(0.6 - 0.2, 0.8 - 0.2, 0.3 - 0.1) \\ &= 0.5 + 0.6 = 1.1 \text{ Mbytes.} \end{aligned}$$

This memory requirement is lower than the available memory, meaning that, provided the tasks are preempted based on their memory usage, all three tasks can be 25 executed concurrently.

Accordingly, the task manager 503 invokes “memory-based preemption mode” by instructing the tasks to transmit deschedule instructions to the scheduler 501 at their designated preemption points $MP_{i,j}$. In this mode, the scheduler 501 allows each task to run non-preemptively from one preemption point to the next, with the constraint that, at 30 any point in time, at most one task at a time can be at a point other than one of its preemption points. Assuming that the newly arrived task starts at a preemption point, the

scheduler 501 ensures that this condition holds for the currently running tasks, thereby constraining all but one task to be at a preemption point.

Thus, in one known memory-based preemption mode, the scheduler 501 is only allowed to preempt tasks at their memory preemption points (i.e. in response to a
5 deschedule request from the task at their memory-based preemption points).

When one of the tasks has terminated, the terminating task informs the task manager 503 that it is terminating, causing the task manager 503 to evaluate Equation 1 and if the worst case memory usage (taking into account removal of this task) is lower than that available to the scheduler 501, the task manager 503 can cancel memory-based
10 preemption, which has the benefit of enabling the system to react faster to external events (since the processor is no longer "blocked" for the duration of the sub-jobs). In general, termination of a task is typically caused by its environment, e.g. a switch of channel by the user or a change in the data stream of the encoding applied (requiring another kind of decoding), meaning that the task manager 503 and/or scheduler 501 should be aware of the
15 termination of a task and probably even instruct the task to terminate.

It should be noted that, when invoked, memory-based preemption constraints are obligatory.

The tasks have been described as software tasks, but a task can also be implemented in hardware. Typically, a hardware device (behaving as a hardware task) is
20 controlled by a software task, which allocates the (worst-case) memory required by the hardware device, and subsequently instructs the hardware task to run. When the hardware task completes, it informs the software task, which subsequently de-allocates the memory. Hence, by having a controlling software task, hardware tasks can simply be dealt with as described above.

25 This approach applies to a single processor system and is not optimized for a multi-processor system and thus a multi-processor optimization is needed.

The present invention provides optimizations of the single processor approach to memory based resource management described above, for multi-processor systems. Consider a set Γ of n tasks τ_i ($1 \leq i \leq n$) and a set P of p processors π_k ($1 \leq k \leq p$), where n
30 is typically much larger than p and assume a fixed-priority preemptive scheduling (FPPS) scheme for the non-constrained mode, and a fixed-priority scheduling scheme with deferred preemption (FPDP) scheme for the memory-based preemption mode.

There are three alternative preferred embodiments, depending on the allocation of tasks to processors:

1. Fixed Allocation: every task τ_i is allocated to a particular processor π_k , i.e. task τ_i will exclusively execute on processor π_k . This embodiment is preferred when processors are dedicated, i.e., where each processor differs essentially from every other processor.
2. Variable Allocation: every task τ_i may execute on every processor π_k . At run-time the scheduler determines which processor executes which task. A task may be preempted while running on one processor, and later continue on another. This embodiment is preferred when all the processors are identical.
3. Mixed Allocation: every task τ_i is allocated to a subset of processors. This is a natural approach when the set of processors can be divided into subsets in which the processors are identical.

The foregoing and other features and advantages of the invention will be apparent from the following, more detailed description of preferred embodiments as illustrated in the accompanying drawings in which reference characters refer to the same parts throughout the various views.

FIG. 1 illustrates a schematic diagram of components of a task interface according to an embodiment of the present invention;

FIG. 2 illustrates a schematic diagram of an example of a digital television system in which an embodiment of the present invention is operative;

FIG. 3 illustrates a schematic diagram of the relationships between components of the task interface illustrated in FIG. 1 for a single processor.

FIG. 4A illustrates components constituting the set-top box of FIG. 2, for a single processor system.

FIG. 4B illustrates components constituting the set-top box of FIG. 2 for a multiprocessor system.

FIG. 5 illustrates components of the processor of the set-top box illustrated in FIG. 2 and FIG. 4A.

It is to be understood by persons of ordinary skill in the art that the following descriptions are provided for purposes of illustration and not for limitation. An artisan understands that there are many variations that lie within the spirit of the invention and the scope of the appended claims. Unnecessary detail of known functions and operations may
5 be omitted from the current description so as not to obscure the present invention.

High volume electronic (HVE) consumer systems, such as digital TV sets, digitally improved analog TV sets and set-top boxes (STBs) must provide real-time services while remaining cost-effective and robust. Consumer products, by their nature, are heavily resource constrained. As a consequence, the available resources have to be used very
10 efficiently, while preserving typical qualities of HVE consumer systems, such as robustness, and meeting stringent timing requirements. Concerning robustness, no one expects, for example, a TV set to fail with the message "please reboot the system".

Significant parts of the media processing in HVE consumer systems are implemented in on-board software that handles multiple concurrent streams of data, and in
15 particular must very efficiently manage system resources, such as main memory, in a multi-tasking environment. Consider a set-top box as an example of an HVE consumer system requiring real-time resource management. Conventionally, as illustrated in FIG. 2, a set-top box 200 receives input for television 201 from a content provider 203 (a server or cable) and from a user interface 205. The user interface 205 comprises a remote control
20 interface for receiving signals from a user-controlled remote device 202, e.g., a handheld infrared remote transmitter. The set-top box 200 receives at least one data stream from at least one of an antenna and a cable television outlet, and performs at least one of processing the data stream or forwarding the data stream to television 201. A user views the at least one data stream displayed on television 201 and via user interface 205, makes
25 selections based on what is being displayed. The set-top box 200 processes the user selection input and based on this input may transmit to the content provider 203 the user input, along with other information identifying the set-top 200 and its capabilities.

FIG. 4B illustrates a simplified block diagram of an exemplary system 450 of a typical set-top box 200 that may include a plurality of processors 460.1 - 460.3 managed
30 by a control and allocation logic module which allocates tasks to processors 460.1 - 460.3 and controls the overall operation of set-top box 200. The control and allocation logic module 451 comprises a data receiver 452 for receiving put data from the network 407 and from storage 406, an evaluator 453 for evaluating memory usage, an allocator 454 for

allocating tasks to processors, a selector 455 for selecting tasks to initiate and terminate execution thereof and a scheduler 501 for scheduling tasks for execution and descheduling executing tasks. The control and allocation logic module 451 is coupled to a television tuner 403, a memory 405, a long term storage device 406, a communication interface 407, and a remote interface 409. The television tuner 403 receives television signals over transmission line 411 and these signals may originate from at least one of an antenna (not shown) and a cable television outlet (not shown). The control & allocation logic module 451 manages the user interface 205, providing data, audio and video output to the television 201 via line 413. The remote interface 409 receives signals from the remote control via the wireless connection 415. The communication interface 407 interfaces between the set-top box 200 and at least one remote processing system, such a Web server, via data path 417. The communication interface 417 is at least one of a telephone modem, an Integrate Services Digital Network (ISDN) adapter, a Digital Subscriber Line (xDSL), a cable television modem, and any other suitable data communication device. The exemplary system 450 of FIG. 4B is for descriptive purposes only. Although the description may refer to terms commonly used in describing particular set-top boxes 200, the description and concepts equally apply to other control processors, including systems having architectures dissimilar to that shown in FIG. 4B.

The control and allocation logic module 451, in a preferred embodiment, is configured to allocate a plurality of real-time tasks relating to the control of the set-top box 200 to a plurality of multi-processors 460.1 - 460.3 that share a memory 405. These real-time tasks include changing channels, selection of a menu option displayed on the user interface 205, decoding incoming data streams, recording incoming data streams using the long term storage device 406 and replaying them, etc. The operation of the set-top box is determined by these real-time control tasks based on characteristics of the set-top box 100, incoming video signals via line 411, user inputs via user interface 205, and any other ancillary input.

In a preferred embodiment, multi-processors share memory 405. These multi-processors 460.1 - 460.3 may each be CPUs, or co-processors, or other processing devices. In a preferred embodiment, the same task is never executed simultaneously by two (or more) processors. In a preferred embodiment, the control and allocation logic module comprises a single task-manager for the entire set of processors 460.1 - 460.3. This is a centralized approach for discussion purposes only. The present invention is not

constrained to centralized approaches, and decentralized versions may be easily conceived by individuals experienced in the art.

Fixed allocation

Whenever every task τ is allocated to a particular processor π , the set Γ of tasks
 5 may be partitioned into p disjoint sets $\Gamma_1 - \Gamma_p$, where the subset Γ_k is allocated to processor π_k . The maximum amount of memory required by the subset of tasks Γ_k under FPPS as executed by processor π_k is given by (Eq. 1s'), which is a variant of (Eq. 1). The term n_k in (Eq. 1s') denotes the number (i.e. the cardinality of the subset Γ_k) of tasks allocated to processor π_k , and the variable i is assumed to range over the tasks of Γ_k .

$$M_k^P = \sum_{i=1}^{n_k} \max_{j=1}^{m(i)} MI_{i,j} \quad (\text{Eq. 1s'})$$

The total memory requirements M^P of the entire set Γ of tasks is determined by adding the results found for each subset of tasks; see (Eq. 1m^{fix}).

$$M^P = \sum_{k=1}^p M_k^P \quad (\text{Eq. 1m}^{\text{fix}})$$

When M^P does not exceed the available memory M_{sys} , there is no need to constrain the scheduling of the tasks on any of the processors. When M^P does exceed M_{sys} , we may
 20 constrain the scheduling of one or more tasks on one or more processors. The effect of constraining the scheduling of all tasks on a single processor can be determined using Eq. 2s'), which is a variant of (Eq. 2s).

$$M_k^D = \sum_{i=1}^{n_k} \max_{j=1}^{m(i)} MP_{i,j} + \max_{1 \leq i \leq n_k} (\max_{j=1}^{m(i)} MI_{i,j} - \max_{j=1}^{m(i)} MP_{i,j}) \quad (\text{Eq. 2s'})$$

The effect of constraining the scheduling of all tasks on all processors can be determined by (Eq. 2m^{fix}), where the total memory requirements M^D is the sum of the memory requirements M_k^D of each set Γ_k .

$$M^D = \sum_{k=1}^p M_k^D \quad (\text{Eq. 2m}^{\text{fix}})$$

Clearly, there are many intermediate alternative embodiments, such as constraining all the tasks on just a subset of the processors, and constraining only a subset of tasks on a subset of the processors.

5

Variable allocation

Alternatively, every task may be executed on every processor 460.1 - 460.3, and scheduling of the tasks is performed by the control and allocation logic module 451. The maximum memory requirements M^P for the non-constrained mode remains the same, i.e.

10 can be determined using (Eq. 1s)

$$M^P = \sum_{i=1}^n \max_{j=1}^{m(i)} MI_{i,j} \quad (\text{Eq. 1s})$$

The constrained mode (i.e. all tasks are only preempted at their preemption points) 15 requires a variant of (Eq. 2s)

$$\begin{aligned} M^D &= \sum_{i=1}^n \max_{j=1}^{m(i)} MP_{i,j} + \max_{i=1}^3 (\max_{j=1}^{m(i)} MI_{i,j} - \max_{j=1}^{m(i)} MP_{i,j}) \quad (\text{Eq. 2s}) \\ &= 0.2 + 0.2 + 0.1 + \max(0.6 - 0.2, 0.8 - 0.2, 0.3 - 0.1) \\ &= 0.5 + 0.6 = 1.1 \text{ Mbytes.} \end{aligned}$$

20

because p tasks may now run in parallel on p processors. The total memory requirements M^D is now given by (Eq. 2m^{var}).

$$25 \quad M^D = \sum_{i=1}^n \max_{j=1}^{m(i)} MP_{i,j} + \max_{1 \leq i_1 < i_2 < \dots < i_p \leq n} \sum_{k=1}^p (\max_{j=1}^{m(i_k)} MI_{i_k,j} - \max_{j=1}^{m(i_k)} MP_{i_k,j}) \quad (\text{Eq. 2m}^{\text{var}})$$

Note that (Eq. 2m^{var}) assumes $n \geq p$, and that (Eq. 2m^{var}) is also only valid for $n \geq p$. For the special case $n < p$, the equation can be simplified to

$$30 \quad M^D = \sum_{i=1}^n \max_{j=1}^{m(i)} MI_{i,j}$$

Clearly, there are many intermediate embodiments, such as constraining only a subset of tasks.

As an example, consider the case of two processors (i.e. $p = 2$), and three tasks (i.e. $n = 3$). It is assumed that the tasks have the same characteristics (i.e. memory requirements) as described in Tables 1 and 2 above, and their accompanying discussion. The maximum memory requirements M^P for the non-constrained mode remains the same, i.e. can be determined using (Eq. 1), and hence exceeds the available memory M_{sys} . Using (Eq. 2m^{var}) the maximum memory requirements M^D for memory-based preemption is:

$$\begin{aligned} M^D &= \sum_{i=1}^3 \max_{j=1}^{m(i)} MP_{i,j} + \max_{1 \leq i_1 < i_2 \leq 3} \sum_{k=1}^2 (\max_{j=1}^{m(i_k)} MI_{i_k,j} - \max_{j=1}^{m(i_k)} MP_{i_k,j}) \\ &= 0.2 + 0.2 + 0.1 + \max(0.4 + 0.6, 0.4 + 0.2, 0.6 + 0.2) \\ &= 0.5 + 1.0 = 1.5 \text{ Mbytes} \end{aligned}$$

The memory requirement is lower than the available memory, meaning that provided that the tasks are preempted only at their preemption points, all three tasks can be executed concurrently.

Mixed Allocation

Assume there are s pair-wise disjoint subsets P_1, \dots, P_s of processors. Let every task be allocated to a particular subset P_l of processors, where $1 \leq l \leq s$. The set of tasks may therefore be divided in s pair-wise disjoint subsets $\Gamma_1, \dots, \Gamma_s$ of tasks. For ease of presentation, the tasks in subset Γ_1 are denoted by τ_i , where $1 \leq i \leq n_1$, i.e. the tasks per subset of tasks are numbered. Similar to the variable allocation described above, the maximum memory requirements M_l^P for subset P_l for the non-constrained mode can be determined using (Eq. 1s'').

The maximum amount of memory required by the subset of tasks Γ_1 in the non-constrained mode executed by the subset P_l of processor π_k is given by (Eq. 1s''), which is a variant of (Eq. 1).

$$M_l^P = \sum_{i=1}^{n_l} \max_{j=1}^{m(i)} MI_{i,j} \quad (\text{Eq. 1s''})$$

Note that (Eq. 1s'') is similar, but not identical, to (Eq. 1s'). Whereas l ranges over subsets of processors in the former equation, k ranges over the processors in a subset in the latter equation. The total memory requirements M^P for all s subsets of processors can be found by taking the sum of the requirements for the individual subsets:

5

$$M^P = \sum_{l=1}^s M_l^P \quad (\text{Eq. 1m}^{\text{mix}})$$

When M^P exceeds M_{sys} , the scheduling is constrained of one or more tasks of one or more subsets of tasks executed on their associated subsets of processors. The effect of constraining the scheduling of all n_1 tasks of Γ_1 on a single subset P_1 of p_1 processors can be determined using (Eq. 2m^{mix}), which is similar to (Eq. 2m^{var}):

$$M_1^D = \sum_{i=1}^{n_1} \max_{j=1}^{m(i)} MP_{i,j} + \max_{1 \leq i_1 < i_2 < \dots < i_p \leq n_1} \sum_{k=1}^{p_1} (\max_{j=1}^{m(i_k)} MI_{i_k,j} - \max_{j=1}^{m(i_k)} MP_{i_k,j}) \quad (\text{Eq. 2m}^{\text{mix}})$$

15

Similarly to (Eq. 2m^{var}), (Eq. 2m^{mix}) only holds for $n_1 \geq p_1$.

The effect of constraining the scheduling of all tasks of every subset of tasks on every subset of processors can be determined by (Eq. 2m^{mix'}), where the total memory requirements M^D is the sum of the memory requirements M_l^D of each subset of processor P_l .

20

$$M^D = \sum_{l=1}^s M_l^D \quad (\text{Eq. 2m}^{\text{mix}'})$$

25

As for the fixed allocation and variable allocation cases described above, for mixed allocation there are many intermediate solutions between the non-constrained mode and the memory-based preemption mode in which all tasks are only preempted at their preemption points.

30

While the preferred embodiments of the present invention have been illustrated and described, it will be understood by those skilled in the art that various changes and

modifications may be made, and equivalents may be substituted for elements thereof without departing from the true scope of the present invention. In addition, many modifications may be made to adapt the teaching of the present invention to a particular situation without departing from its central scope. Therefore it is intended that the present

5 invention not be limited to the particular embodiments disclosed as the best mode contemplated for carrying out the present invention, but that the present invention include all embodiments falling within the scope of the appended claims.

REFERENCES

The following references support corresponding references numbers in the text and are hereby included herein by reference as if fully set forth herein:

- [1] R. Gopalakrishnan and G.M. Parulkar, "Bringing Real-Time Scheduling Theory
5 And Practice Closer For Multimedia Computing," In: Proc. ACM Sigmetrics Conf. on
Measurement & modeling of computer systems, pp. 1-12, May 1996.
- [2] S. Lee, C.-G. Lee, M. Lee, S.L. Min, and C.-S. Kim, "Limited Preemptible
Scheduling to Embrace Cache Memory In Real-Time Systems," In: Proc. ACM Sigplan
10 Workshop on Languages, Compilers and Tools for Embedded Systems (LCTES), LNCS-
1474, pp. 51-64, June 1998.
- [3] J. Simonson and J.H. Patel, "Use Of Preferred Preemption Points In Cache-Based
Real-Time Systems", In: Proc IEEE International Computer Performance and
15 Dependability Symposium (IPDS'95), pp. 316-325, April 1995.
- [4] R. J. Bril and D. J. C. Lowet, "A Method For Handling Preemption Points," Philips
Research Laboratories, Eindhoven, The Netherlands, Internal IST/IPA document, 30
September 2002.
20
- [5] R. J. Bril and D. J. C. Lowet, "A Method For Handling Preemption Points –
Remarks –, " Philips Research Laboratories, Eindhoven, The Netherlands, Internal IST/IPA
document, 31 October 2002.
- 25 [6] Clemens Szyperski, Component Software - Beyond Object-oriented Programming,
Addison-Wesley, ISBN 0-201-17888-5, 1997.

CLAIMS:

1. A method of scheduling a plurality of tasks in a data processing system having a plurality of processors, comprising the steps of:
for each task of the plurality, providing suspension data specifying suspension of the task based on memory used thereby;
allocating each of said plurality of tasks to a processor of said plurality of processors; and
each processor performing the steps of:
 - (i) processing one the tasks allocated to the processor,
 - (ii) monitoring for an input indicative of memory used by the task matching the suspension data associated with the task,
 - (iii) suspending said task on the basis of said monitored input, and
 - (iv) processing a different one of the tasks allocated to the processor.
2. The method of claim 1, wherein:
allocation of a task to a processor is based on one of -
 - a. fixed allocation having every task allocated to a particular processor;
 - b. variable allocation having every task allocated to every processor; and
 - c. mixed allocation having every task allocated to a subset of said plurality of processors.
3. The method of claim 2, wherein said input comprises data indicative of a suspension request.
4. The method of claim 3, further comprising the steps of:
receiving first data identifying maximum memory usage associated with the plurality of tasks;
receiving second data identifying memory available for processing the plurality of tasks;
identifying, on the basis of the first and second data, whether there is sufficient memory available to process the tasks; and

wherein, said monitoring and suspending steps are performed only in response to identifying insufficient memory.

5. The method of claim 4, further comprising the steps of:

monitoring termination of tasks; and

in response to a task termination, repeating said step of identifying availability of memory in response to a task terminating.

6. The method of claim 5, in which, in response to identifying sufficient memory to execute the remaining tasks, the monitoring step is deemed unnecessary.

7. The method of claim 6, further comprising the steps of:

receiving first data identifying maximum memory usage associated with the plurality of tasks;

receiving second data identifying memory available for processing the plurality of tasks;

identifying, on the basis of the first and second data, whether there is sufficient memory available to process the tasks; and

wherein, said monitoring and suspending steps are performed only in response to identifying insufficient memory.

8. The method of claim 7, further comprising the steps of:

monitoring termination of tasks; and

in response to a task termination, repeating said step of identifying availability of memory in response to a task terminating.

9. The method of claim 8, in which, in response to identifying sufficient memory to execute the remaining tasks, the monitoring step is deemed unnecessary.

10. The method of claim 2, further comprising the steps of:

receiving first data identifying maximum memory usage associated with the plurality of tasks;

receiving second data identifying memory available for processing the plurality of tasks;

identifying, on the basis of the first and second data, whether there is sufficient memory available to process the tasks; and

wherein, said monitoring and suspending steps are performed only in response to identifying insufficient memory.

11. The method claim 10, further comprising the steps of:

monitoring termination of tasks; and

in response to a task termination, repeating said step of identifying availability of memory in response to a task terminating.

12. The method according to claim 11, in which, in response to identifying sufficient memory to execute the remaining tasks, the monitoring step is deemed unnecessary.

13. A method of scheduling a plurality of tasks in a data processing system having a plurality of processors, comprising the steps of:

for each task of the plurality, providing suspension data specifying suspension of the task based on memory used thereby;

allocating each of said plurality of tasks to a processor of said plurality of processors based on one of -

- fixed allocation having every task allocated to a particular processor,
- variable allocation having every task allocated to every processor, and
- mixed allocation having every task allocated to a subset of said plurality of

processors; and

each processor performing the steps of:

- (i) processing one the tasks allocated to the processor,
- (ii) monitoring for an input indicative of memory used by the task matching the suspension data associated with the task,
- (iii) suspending said task on the basis of said monitored input, and
- (iv) processing a different one of the tasks allocated to the processor.

14. A scheduler for use in a data processing system having a plurality of processors, the data processing system being arranged to execute a plurality of tasks on said plurality of processors and having access to a specified amount of memory for use in executing the tasks, the scheduler comprising:

a data receiver arranged to receive data identifying maximum memory usage associated with a task;

an evaluator arranged to identify, on the basis of the received data, whether there is sufficient memory to execute the tasks;

an allocator arranged to allocate each of said plurality of tasks to a processor of said plurality of processors based on one of -

- a. fixed allocation having every task allocated to a particular processor,
- b. variable allocation having every task allocated to every processor, and
- c. mixed allocation having every task allocated to a subset of said plurality of processors;

a selector arranged to select at least one task for suspension during execution of the task, said suspension coinciding with a specified memory usage by the task; and

a scheduler 501 arranged to initiate execution of said allocated task and suspend said selected task,

wherein, in response to the evaluator identifying that there is insufficient memory to execute the plurality of tasks, the selector selects at least one task for suspension, on the basis of its specified memory usage, and the specified amount of memory available to the data processing system, and the scheduler suspends execution of the at least one selected task in response to the task using the specified memory.

15. A scheduler according to claim 14, wherein the evaluator is arranged to monitor termination of tasks, and in response to a task terminating, to identify whether there is sufficient memory to execute the remaining tasks.

16. A scheduler according to claim 15, wherein, in response to the evaluator identifying sufficient memory to execute the remaining tasks, the selector is arranged to deselect said selected at least one task.

17. A data processing system having a plurality of processors arranged to execute a plurality of tasks, the data processing system including:

memory arranged to hold instructions and data during execution of a task;

receiving means arranged to receive data identifying maximum memory usage associated with a task and data specifying preemptability of the task;

evaluating means arranged to identify, on the basis of the received data, whether there is sufficient memory to execute the tasks;

an allocator arranged to allocate each of said plurality of tasks to a processor of said plurality of processors based on one of -

- fixed allocation having every task allocated to a particular processor,
- variable allocation having every task allocated to every processor, and
- mixed allocation having every task allocated to a subset of said plurality of processors; and

a scheduler arranged to schedule execution of the tasks on the basis of input received from the evaluating means,

wherein, in response to identification of insufficient memory to execute the plurality of tasks, the scheduler is arranged to suspend execution of at least one task in dependence on memory usage by the task.

18. A method of transmitting data to a data processing system having a plurality of processors, the method comprising:

allocating each task of a plurality of tasks to a processor of said plurality of processors based on one of -

- fixed allocation having every task allocated to a particular processor,
- variable allocation having every task allocated to every processor, and
- mixed allocation having every task allocated to a subset of said plurality of processors;

transmitting data for use by the data processing system in processing each task of said plurality of tasks;

transmitting suspension data specifying suspension of each task of said plurality of tasks based on memory usage during processing thereof;

wherein, the data processing system is configured to perform a process comprising:

monitoring for an input indicative of memory usage of each task matching the suspension data associated with the task; and
suspending processing of said each task on the basis of said monitored input.

19. A method according to claim 18, wherein the suspension data includes data identifying maximum memory usage associated with said each task.

20. A method according to claim 18, wherein the suspension data identifies at least one point at which processing of each task can be suspended, based on memory usage of said each task.

21. A method according to claim 20, wherein said each task comprises a plurality of sub-jobs and said data identifying at least one point at which processing of said each task can be suspended corresponds to each such sub-job.

22. A method according to claim 20, wherein the suspension data includes data identifying maximum memory usage associated with said each task.

23. A method according to claim 22, wherein said each task comprises a plurality of sub-jobs and said data identifying at least one point at which processing of said each task can be suspended corresponds to each such sub-job.

24. A method of configuring a plurality of tasks for use in a data processing system having a plurality of processors, the method including associating suspension data with the task, the suspension data specifying suspension of the task based on memory usage associated therewith, wherein the data processing system is arranged to perform a process in respect of a plurality of tasks executing on a plurality of processors, the process comprising:

allocating each task of said plurality of tasks to a processor of said plurality of processors based on one of -

- fixed allocation having every task allocated to a particular processor,
- variable allocation having every task allocated to every processor, and

- mixed allocation having every task allocated to a subset of said plurality of processors;

monitoring for an input indicative of memory usage of the task matching the suspension data associated with the task; and

suspending processing of said task on the basis of said monitored input.

25. A computer program, comprising a set of instructions arranged to cause a processing system to perform the method according to of claim 1.

26. A computer program, comprising a set of instructions arranged to cause a processing system to perform the method according to of claim 2.

ABSTRACT

A method and apparatus is provided for use by a scheduler of a multi-processor data processing system to select task preemption points based on main memory requirements. There are three alternative preferred embodiments, depending on the allocation of tasks to processors: (1) Fixed Allocation: every task is allocated to a particular processor, i.e. each task exclusively executes on a given processor. This embodiment is preferred when processors are dedicated, i.e., where each processor differs essentially from every other processor; (2) Variable Allocation: every task may execute on every processor. At run-time the scheduler determines which processor executes which task. A task may be preempted while running on one processor, and later continue on another. This embodiment is preferred when all the processors are identical; and (3) Mixed Allocation: every task is allocated to a subset of processors. This is a natural approach when the set of processors can be divided into subsets in which the processors are identical.

IN THE UNITED STATES PATENT AND TRADEMARK OFFICE

In re Application of
REINDER J. BRIL ET AL

Atty. Docket
US 040013

Serial No.:

Filed: CONCURRENTLY

Title: RESOURCE MANAGEMENT IN A MULTI-PROCESSOR SYSTEM

Commissioner for Patents, Alexandria, VA 22313-1450

APPOINTMENT OF ASSOCIATES

Sir:

The undersigned Attorney of Record hereby revokes all prior appointments (if any) of Associate Attorney(s) or Agent(s) in the above-captioned case and appoints:

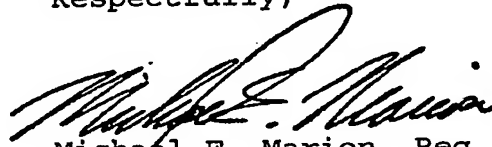
DANIEL J. PIOTROWSKI

(Registration No. 42,079)

c/o U.S. PHILIPS CORPORATION, Intellectual Property Department, 580 White Plains Road, Tarrytown, New York 10591, his Associate Attorney(s)/Agent(s) with all the usual powers to prosecute the above-identified application and any division or continuation thereof, to make alterations and amendments therein, and to transact all business in the Patent and Trademark Office connected therewith.

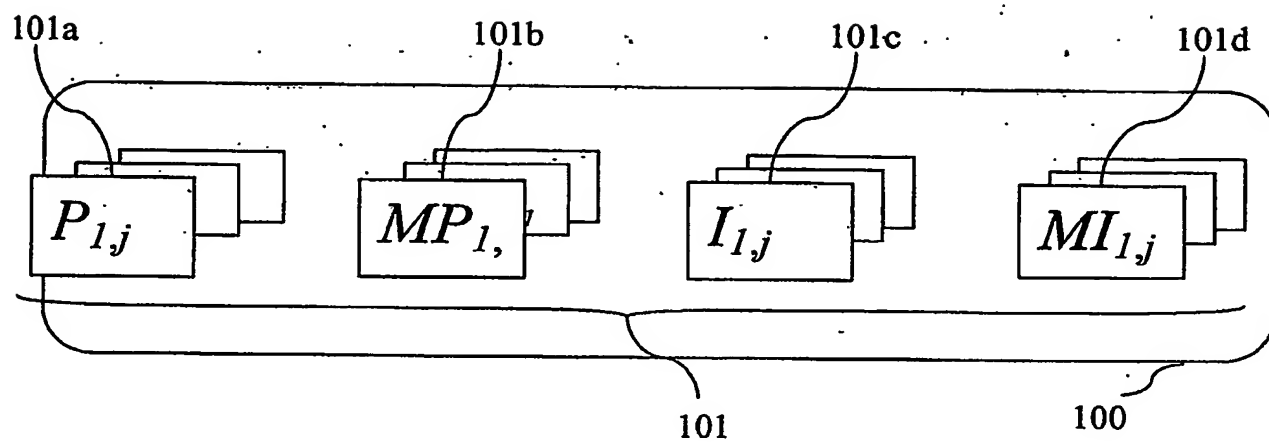
ALL CORRESPONDENCE CONCERNING THIS APPLICATION AND THE LETTERS PATENT WHEN GRANTED SHOULD BE ADDRESSED TO THE UNDERSIGNED ATTORNEY OF RECORD.

Respectfully,



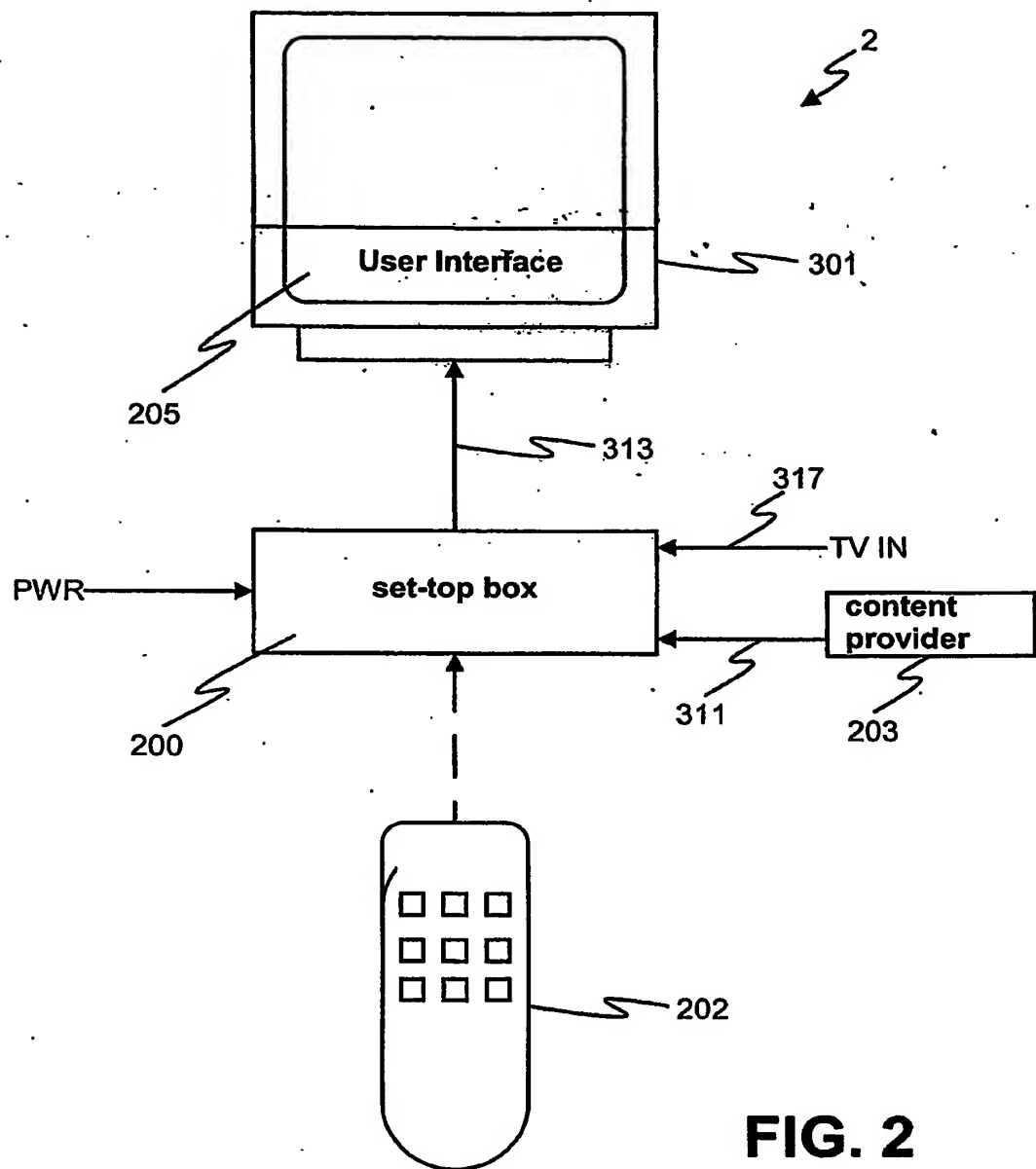
Michael E. Marion, Reg. 32,266
Attorney of Record

Dated at Tarrytown, New York
this 14th day of May, 2003.



Interface for task τ_1

FIG. 1



Task τ_i

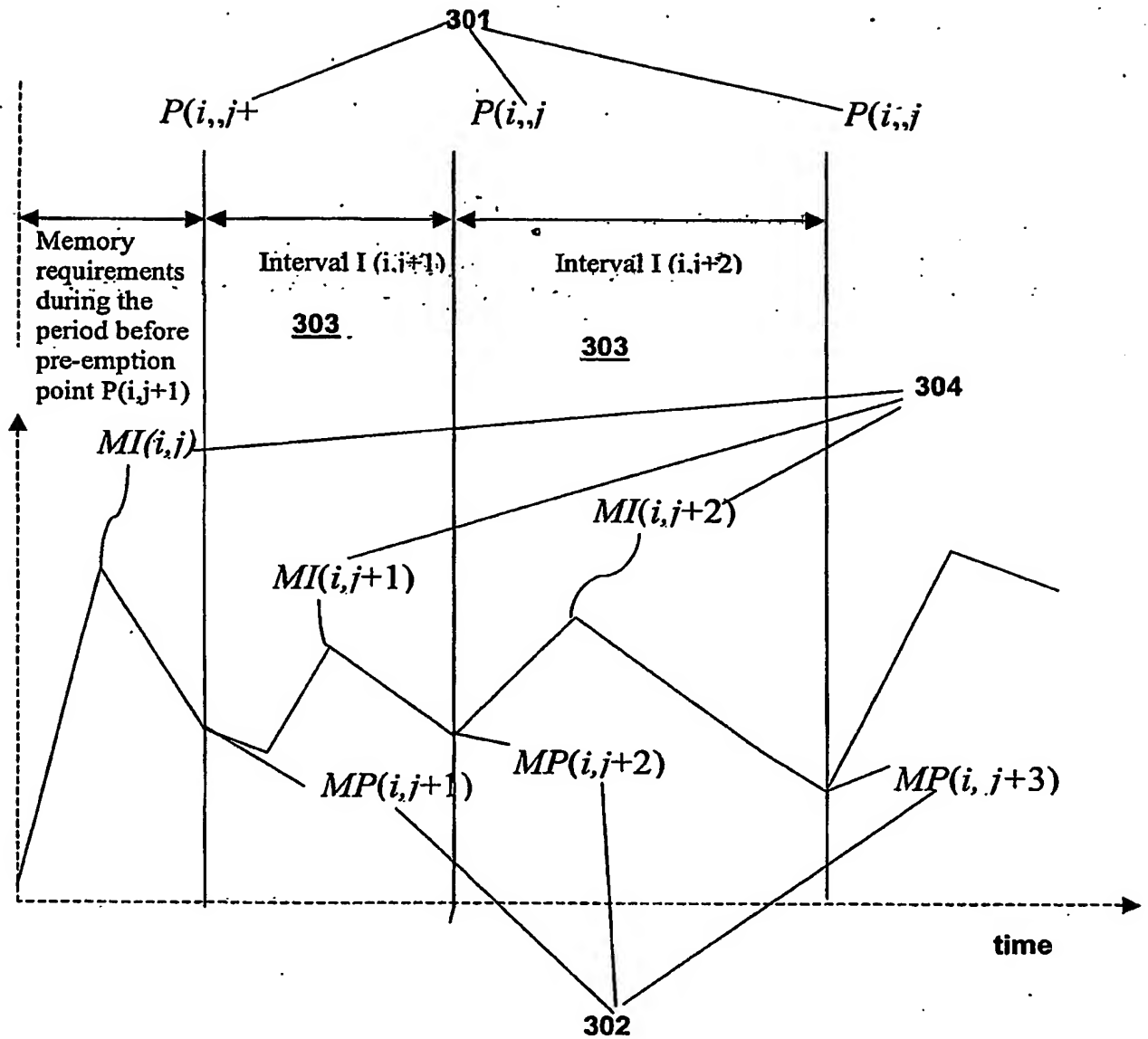


FIG. 3

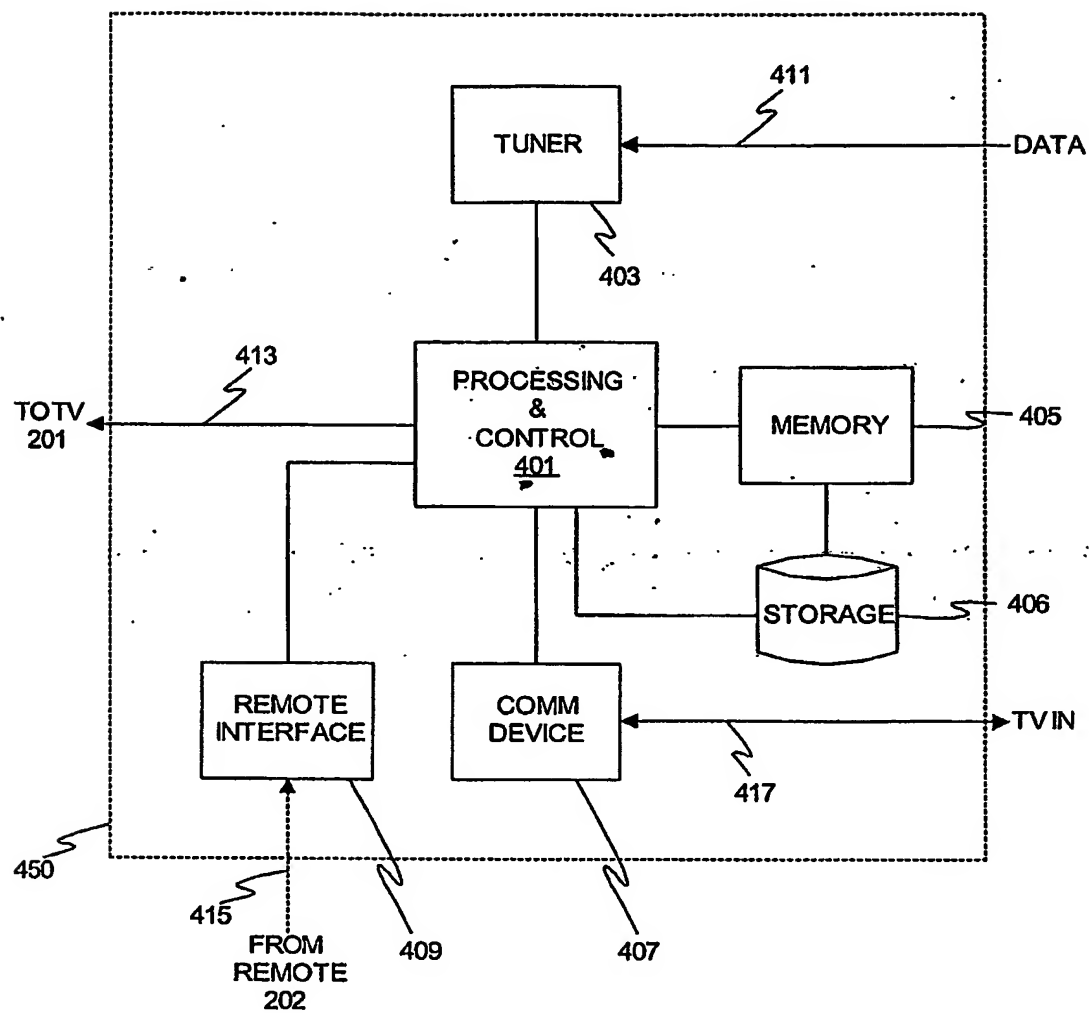


FIG. 4A

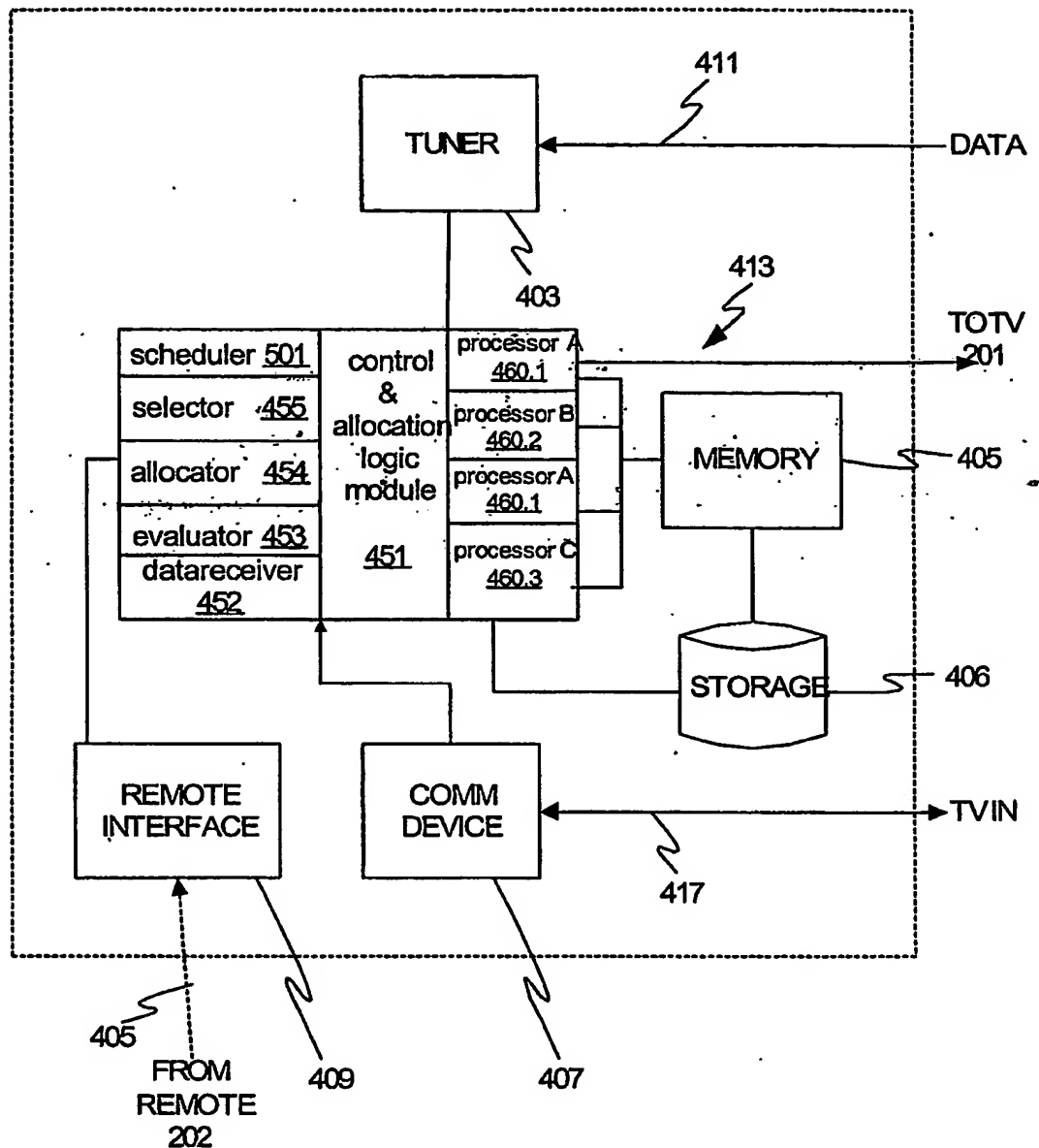


FIG. 4B

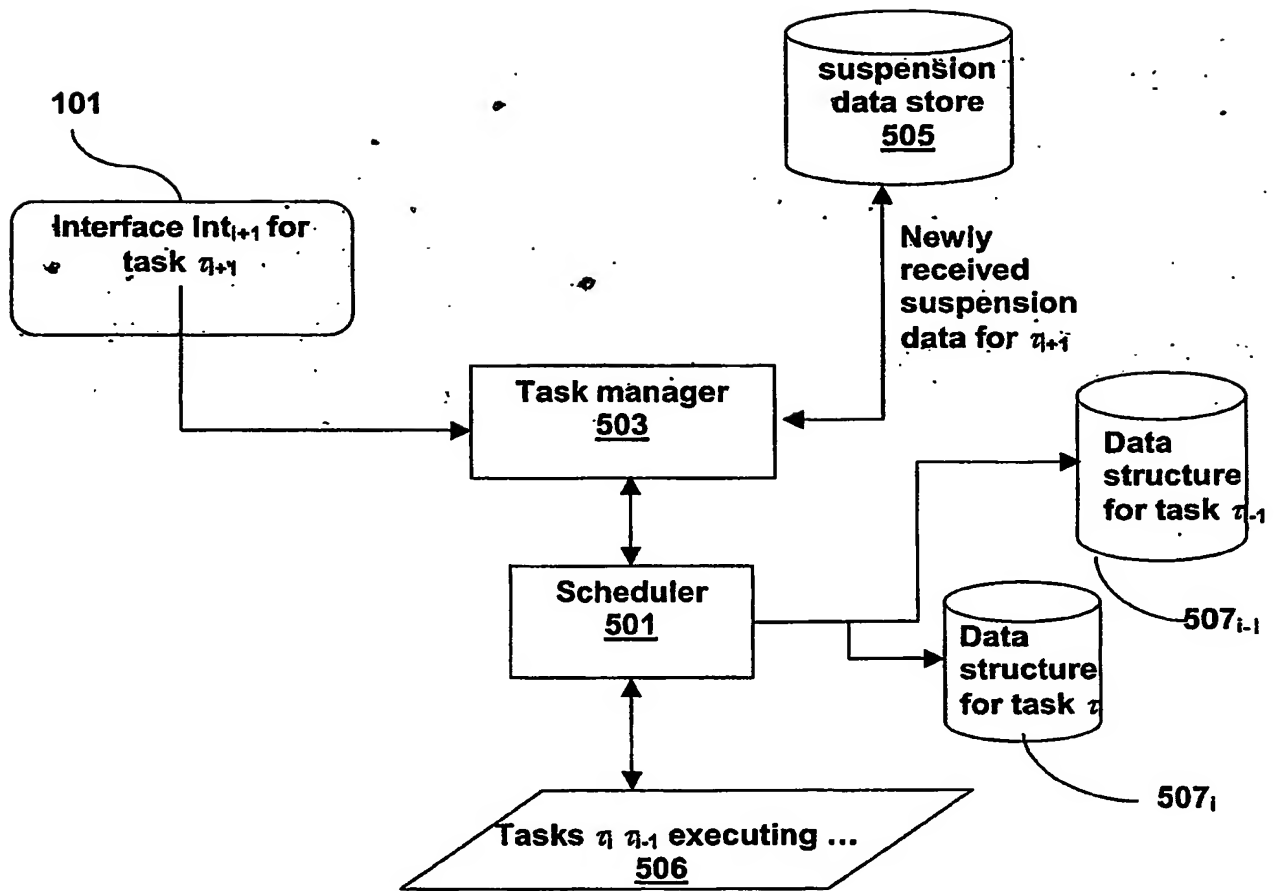


FIG. 5

**This Page is Inserted by IFW Indexing and Scanning
Operations and is not part of the Official Record**

BEST AVAILABLE IMAGES

Defective images within this document are accurate representations of the original documents submitted by the applicant.

Defects in the images include but are not limited to the items checked:

- ☐ **BLACK BORDERS**
- ☐ **IMAGE CUT OFF AT TOP, BOTTOM OR SIDES**
- ☐ **FADED TEXT OR DRAWING**
- ☐ **BLURRED OR ILLEGIBLE TEXT OR DRAWING**
- ☐ **SKEWED/SLANTED IMAGES**
- ☐ **COLOR OR BLACK AND WHITE PHOTOGRAPHS**
- ☐ **GRAY SCALE DOCUMENTS**
- ☒ **LINES OR MARKS ON ORIGINAL DOCUMENT**
- ☐ **REFERENCE(S) OR EXHIBIT(S) SUBMITTED ARE POOR QUALITY**
- ☐ **OTHER: _____**

IMAGES ARE BEST AVAILABLE COPY.

As rescanning these documents will not correct the image problems checked, please do not report these problems to the IFW Image Problem Mailbox.